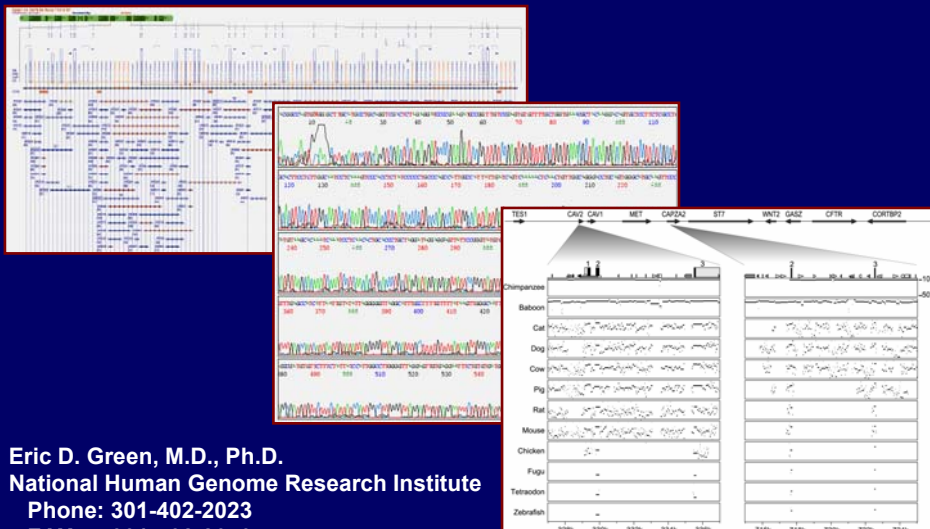
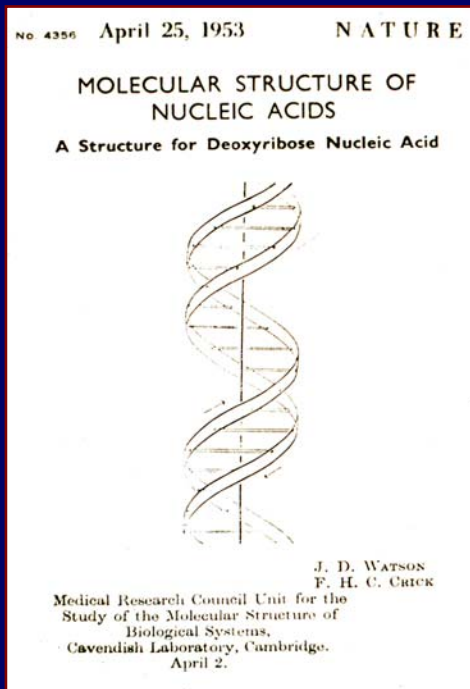
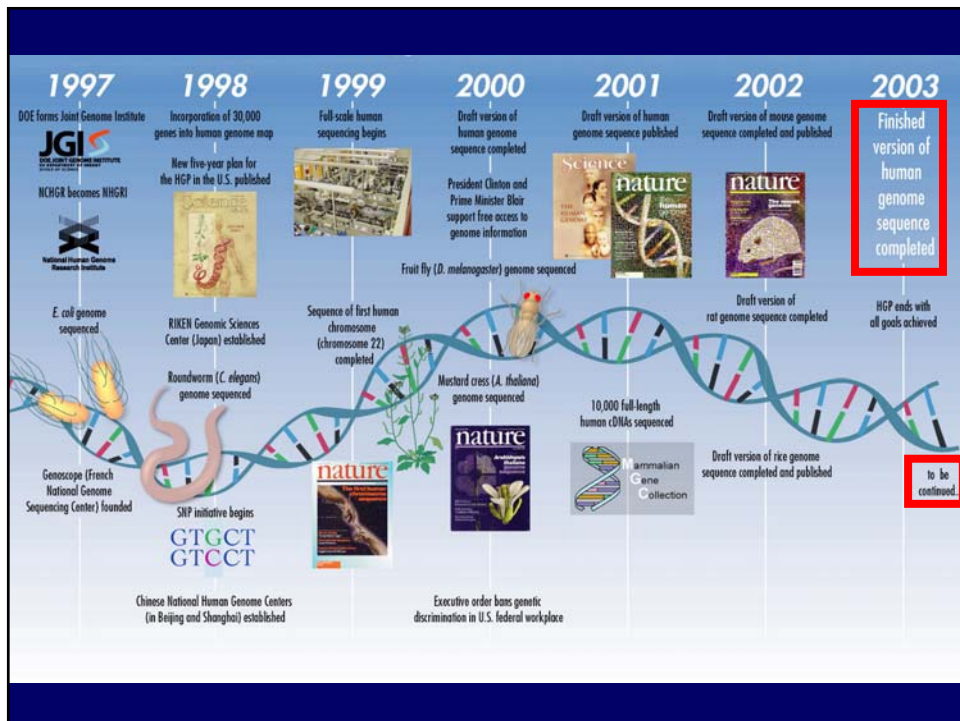
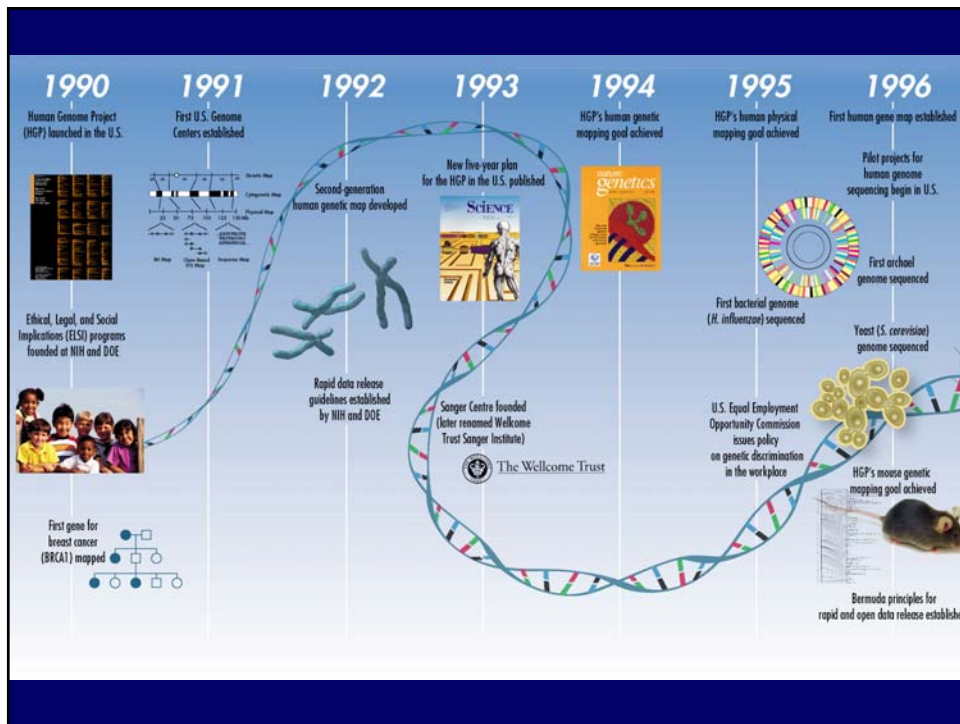


## Techniques for Genome Mapping & Sequencing



Eric D. Green, M.D., Ph.D.  
 National Human Genome Research Institute  
 Phone: 301-402-2023  
 FAX: 301-402-2040  
 E-Mail: egreen@nhgri.nih.gov



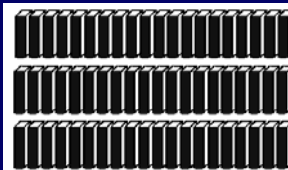


## Outline

- I. Fundamentals of Genome Mapping
- II. Fundamentals of Genome Sequencing
- III. Mapping & Sequencing in the Human Genome Project... and Beyond
- IV. Comparative Sequencing

## Genome Sizes

Human Genome  
Mouse Genome



~3,000,000,000 bp

Fruit Fly Genome



~160,000,000 bp

Nematode Genome



~100,000,000 bp

Yeast Genome

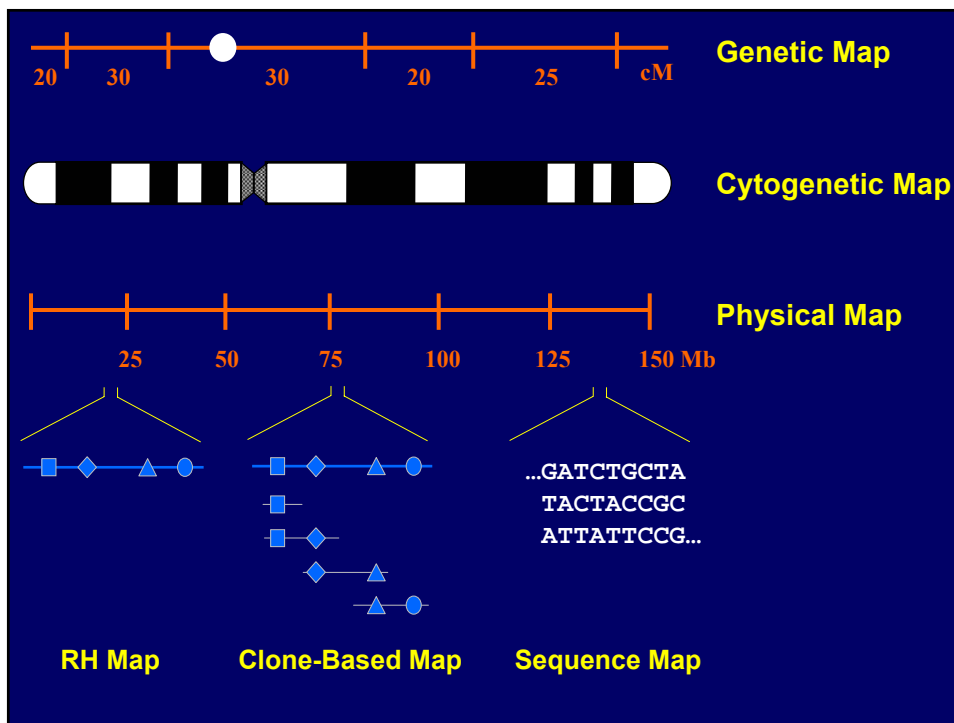
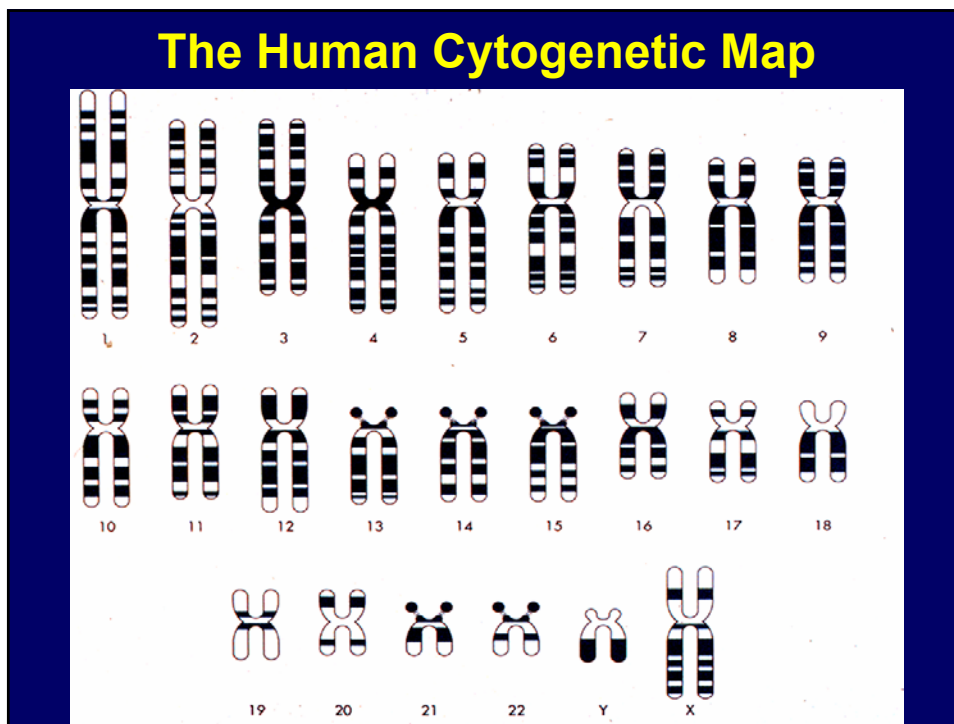


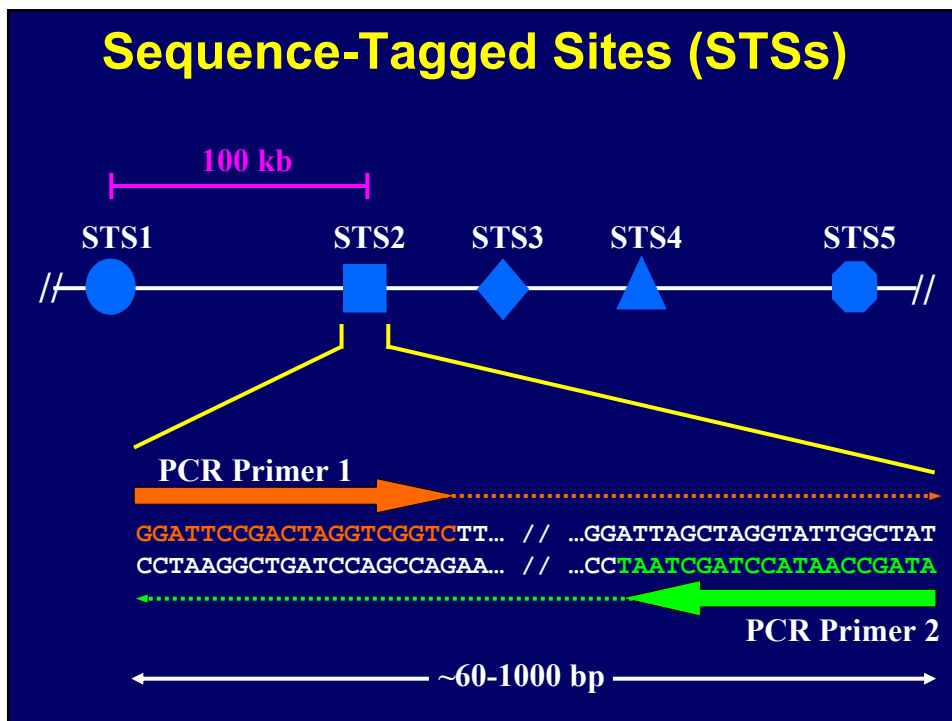
~15,000,000 bp

*E. coli* Genome



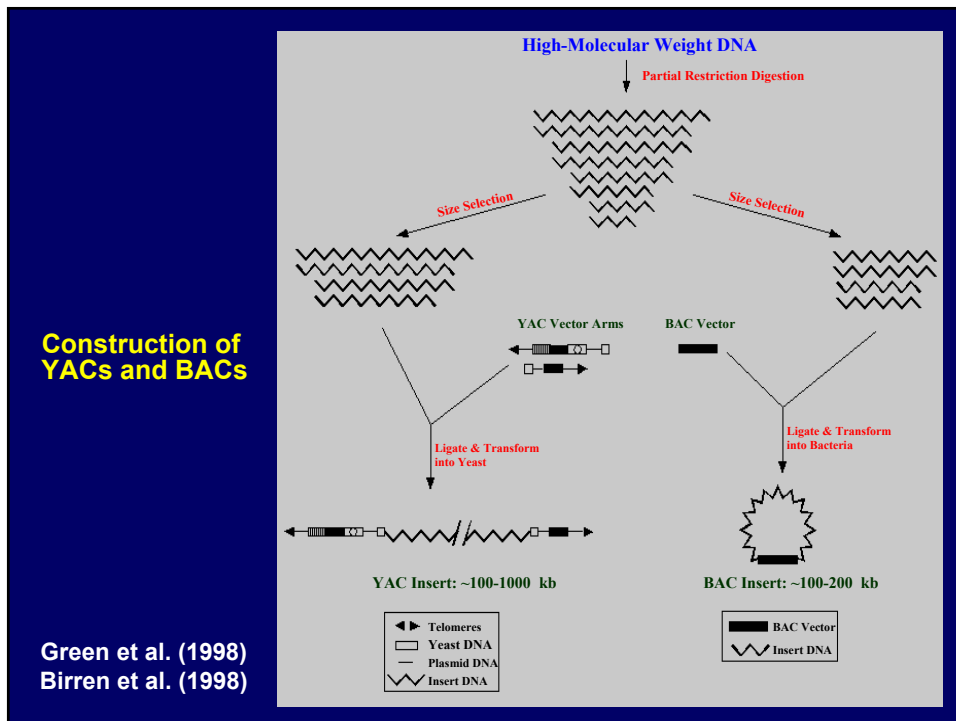
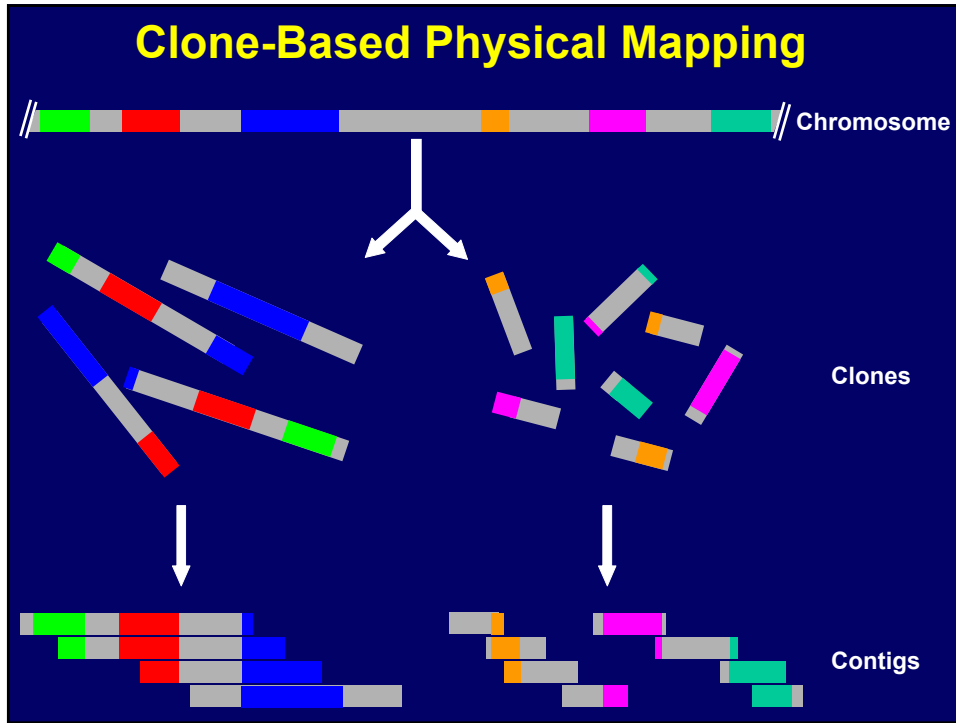
~5,000,000 bp

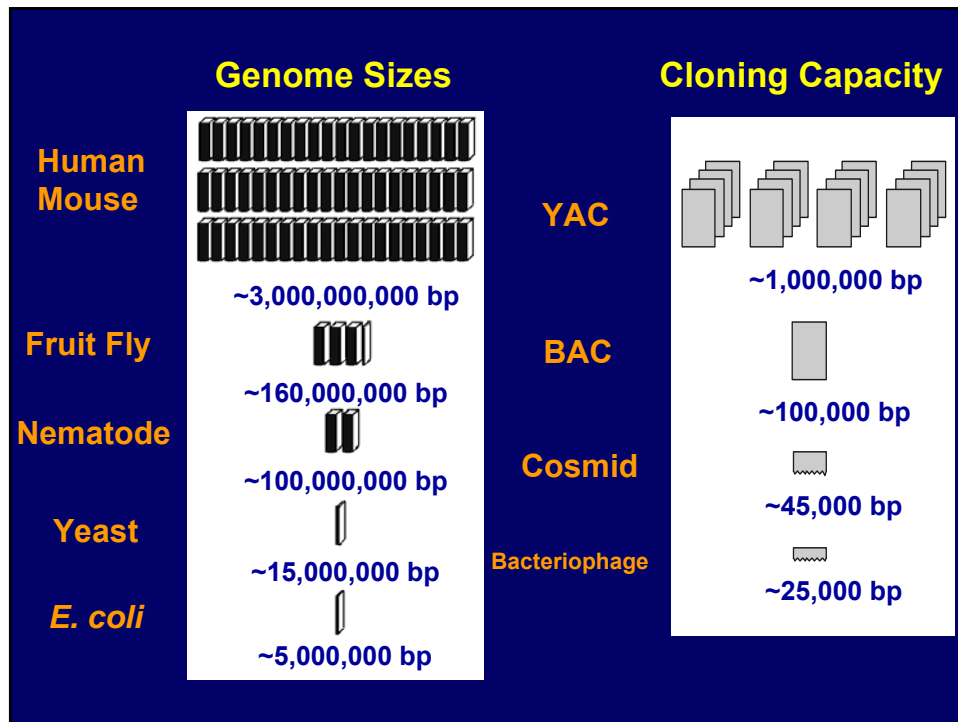




## Physical Mapping: General Principles

- Importance of Physical Maps:
  - Localization and Isolation of Genes (e.g., Positional Cloning)
  - Study of Genome Organization and Evolution
  - Framework for Genome Sequencing
- Physical Mapping Involves Ordering Clones and/or Landmarks
- General Types of Physical Maps:
  - Landmark Only (e.g., Radiation Hybrid Maps)
  - Clone-Based
  - Sequence





## Bacterial Artificial Chromosomes (BACs)

- Bacterial-Based Cloning System Developed by Shizuya et al. (1992)
- Based on the *E. coli* F Factor (Fertility Plasmid): Replication Control
- Cloned Inserts: 100-200 kb, Circular DNA
- Low Copy Number
  - Low Yields of DNA by Standard Methods
  - Reasonably Stable
- Relatively Non-Chimeric
- BAC Libraries from Many Different Species now Available (e.g., [www.chori.org/bacpac](http://www.chori.org/bacpac))
- See Birren et al. (1998)

**Genome**  
(~3000 Mb)

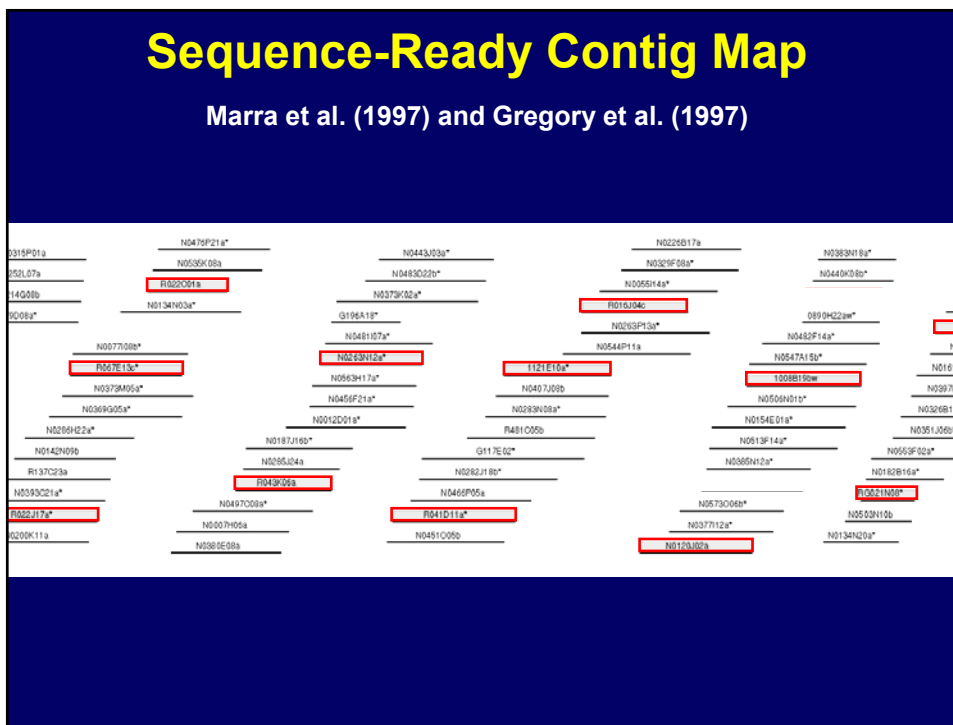
**Chromosome**  
(~130 Mb)

**YAC**  
(~0.5-1.0 Mb)

**BAC**  
(~0.1-0.2 Mb)

g	g	g	g	GATCCTCTAGAATCTC
g	g	g	g	GAGATCTCTAGAGATC
g	g	g	g	GTTGGGAACTCTGTGAA
T	T	T	T	TGTGACTAGCCACAGT
T	T	T	T	TGTGACTAGCCACAGT
T	T	T	T	TACTTGTGAGAGATGT
A	A	A	A	ATGATGCACCTGACCC
g	g	g	g	GGTTCCTACTCTCAAC
g	g	g	g	GACTCAGCTCCACTCTCA
C	C	C	C	CCGGTTAGACATACAT
g	g	g	g	GAGGCCACCCGCCCT
g	g	g	g	GTGCACCTCCACCACC

g	g	g	g	GATCCTCTAGAATCTC
g	g	g	g	GAGATCTCTAGAGATC
g	g	g	g	GTTGGGAACTCTGTGAA
T	T	T	T	TGTGACTAGCCACAGT
T	T	T	T	TGTGACTAGCCACAGT
T	T	T	T	TACTTGTGAGAGATGT
A	A	A	A	ATGATGCACCTGACCC
g	g	g	g	GGTTCCTACTCTCAAC
g	g	g	g	GACTCAGCTCCACTCTCA
C	C	C	C	CCGGTTAGACATACAT
g	g	g	g	GAGGCCACCCGCCCT
g	g	g	g	GTGCACCTCCACCACC

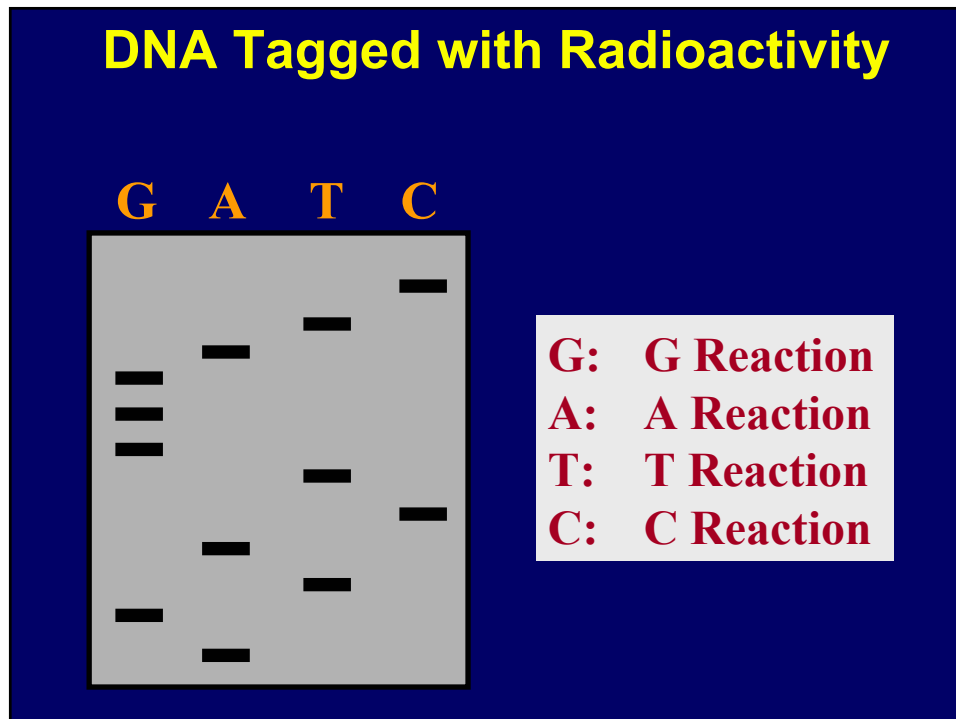
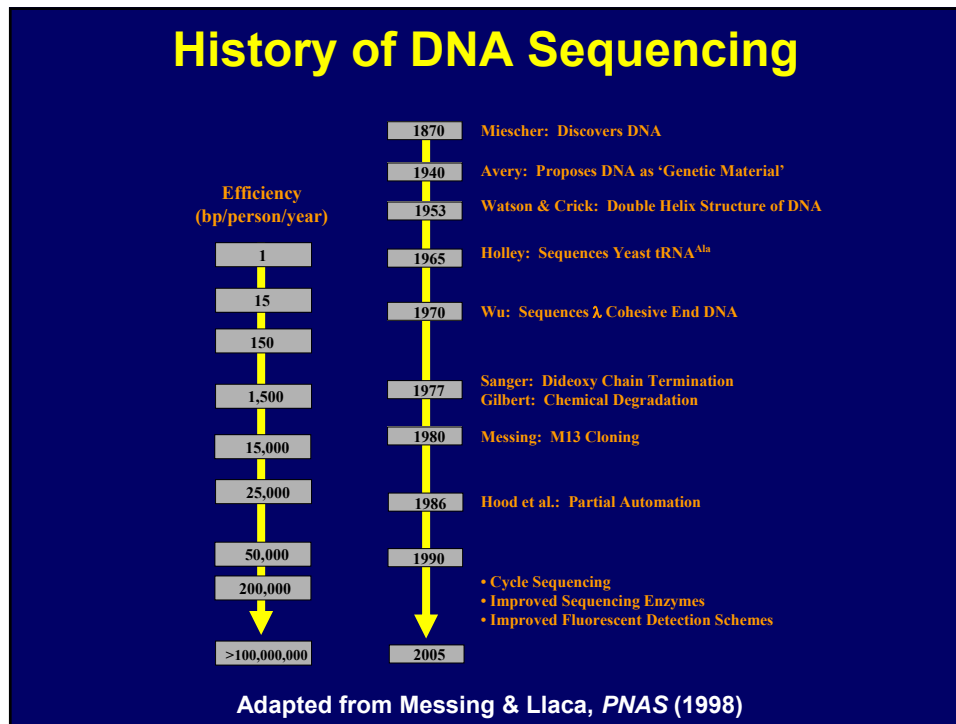




## **Physical Mapping: Future Prospects**

- **Strategies for Physical Mapping are Radically Changing in the Sequence-Based Era**
- **Will Now See a Closer Interplay of Mapping and Sequencing in the Exploration of New Genomes**
- **Construction of New BAC Libraries will Allow Physical Mapping Studies of More Species' Genomes**

## **DNA Sequencing**

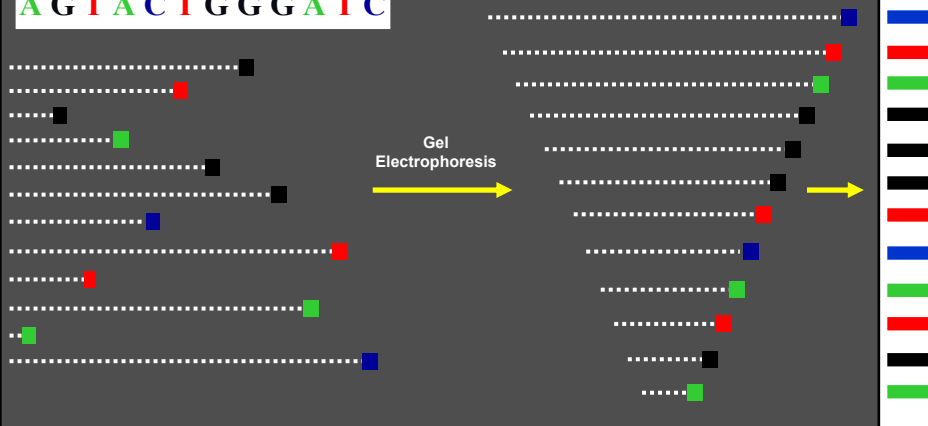


## Radioactive Sequencing

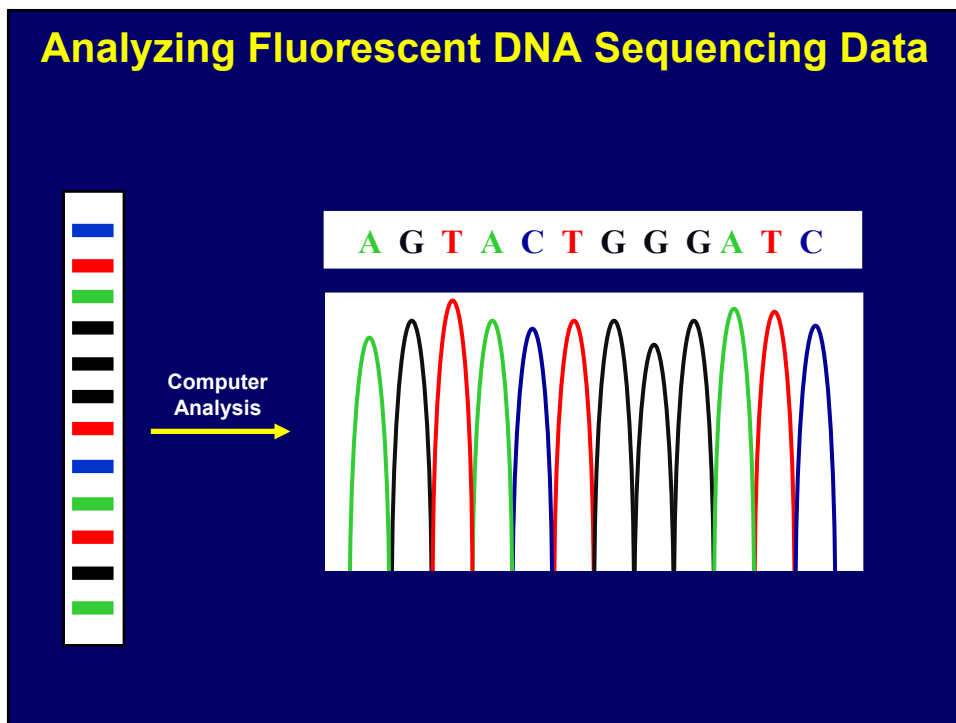
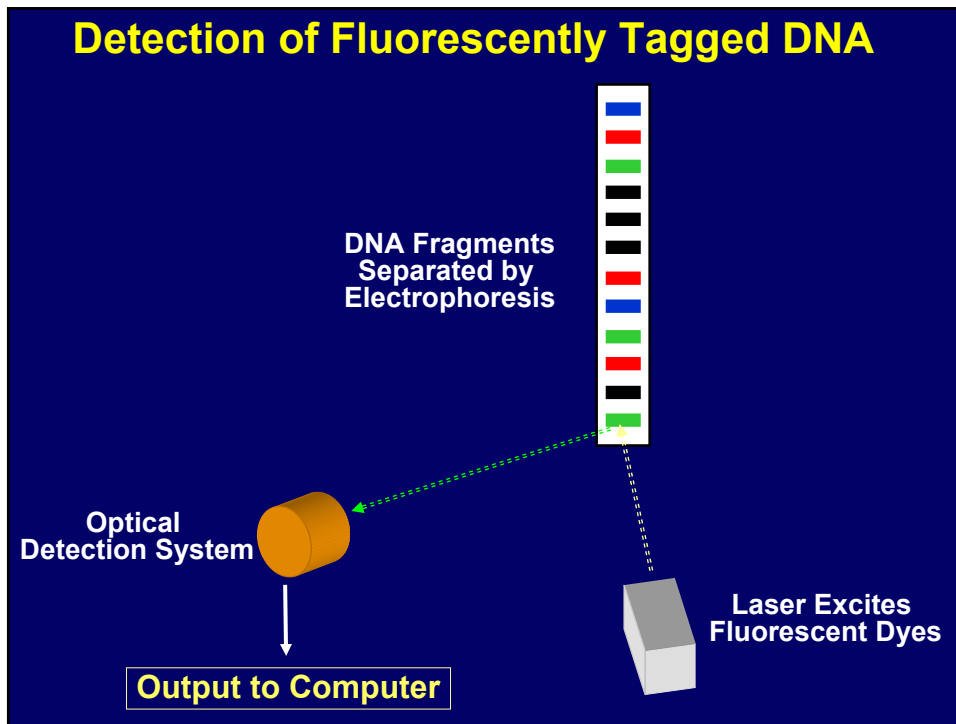


## Fluorescent DNA Sequencing

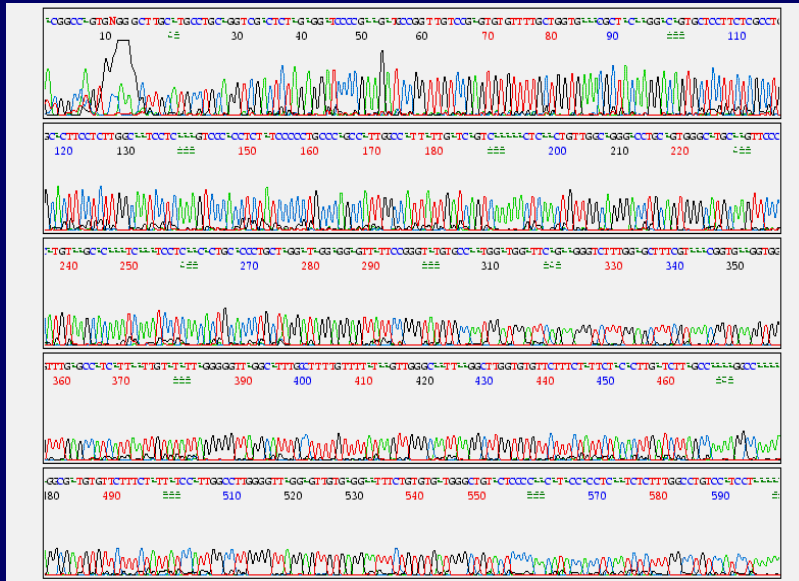
A G T A C T G G G A T C



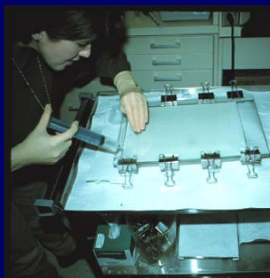
Wilson & Mardis (1997)



## Fluorescent DNA Sequencing Results



## Slab Gel-Based DNA Sequencing Instruments



## Capillary-Based DNA Sequencing Instruments



## Large-Scale cDNA Sequencing

- ESTs: Expressed-Sequence Tags
- SAGE: Serial Analysis of Gene Expression
- Full-Insert (Full-Length) cDNA Sequencing



[mgc.nci.nih.gov](http://mgc.nci.nih.gov)

# Large-Scale Genomic Sequencing



# Shotgun Sequencing

Wilson & Mardis (1997)  
Green (2001)

## Subclone Construction

```
GATGCTTAGAATCTC
GAGTCTTGGAGTCTC
CTGGGAACCTGTGTA
TGTGACTAGCAGAGT
TACTGTGAGAGATAT
ATGATGCACTTACCC
GGTTTACTCTCAGC
GACTCACTCCAGCTCA
GAGGCCACCCGCCCT
CTGACACTTCAGACC
GATTATACCAFTTA
ATCTTAGGATTGACA
```

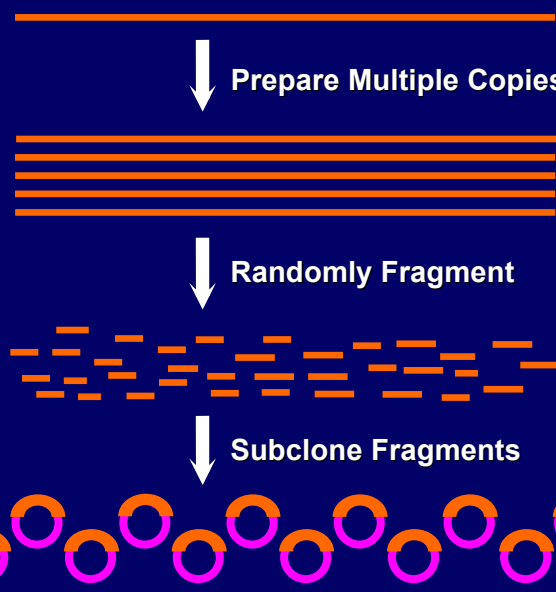
BAC DNA

Prepare Multiple Copies

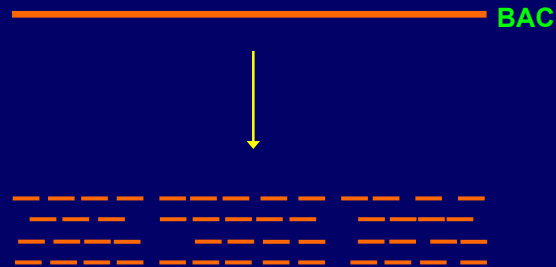
Randomly Fragment

Subclone Fragments

```
GA GA GA GATGCTTAGAATCTC
GA GA GA GAGTCTTGGAGTCTC
GT GT GT CTGGGAACCTGTGTA
TG TG TG TGTGACTAGCAGAGT
TA TA TA TACTGTGAGAGATAT
AT AT AT ATGATGCACTTACCC
GA GA GA GGTTTACTCTCAGC
GA GA GA GACTCACTCCAGCTCA
GA GA GA GAGGCCACCCGCCCT
GT GT GT CTGACACTTCAGACC
GA GA GA GATTATACCAFTTA
AT AT AT ATCTTAGGATTGACA
```



## Shotgun Sequencing Strategy



## Poisson Calculations

The sequencing strategy for the shotgun approach follows the Lander and Waterman application of the Poisson distribution

The probability a base is not sequenced is given by:

$$P_0 = e^{-c}$$

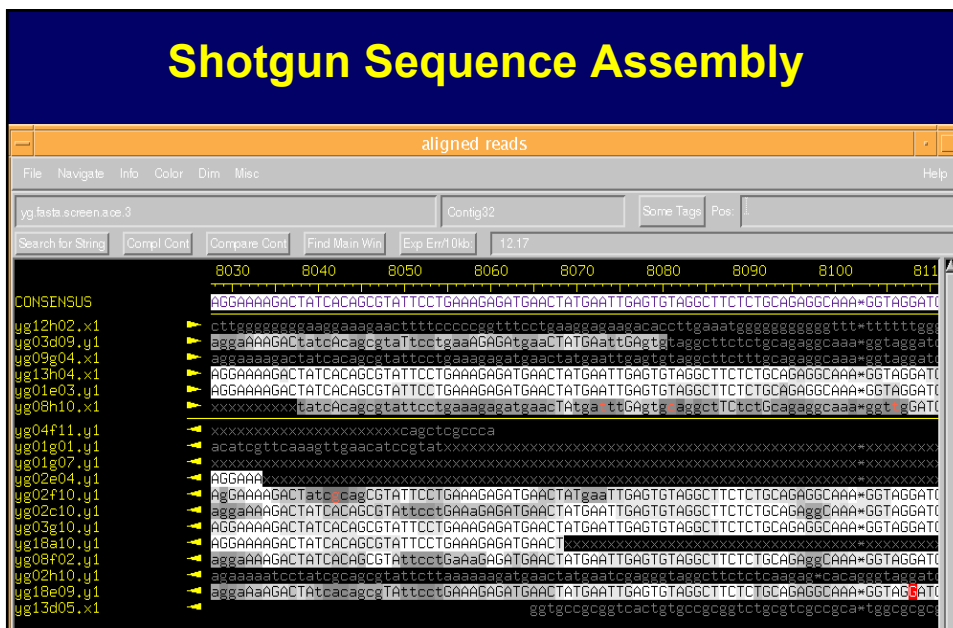
Where:

- $c$  = fold sequence coverage ( $c = LN/G$ ),
- $LN$  = # bases sequenced, i.e.  $L$  = average sequencing read length and  $N$  = # reads
- $G$  = target sequence length
- $e = 2.718$  ( $e = 2.718281828459$ )

Fold Coverage	$P_0 = e^{-c}$	% not sequenced	% sequenced
1	0.37	37%	63%
2	0.135	13.5%	87.5%
3	0.05	5%	95%
4	0.018	1.8%	98.2%
5	0.0067	0.6%	99.4%
6	0.0025	0.25%	99.75%
7	0.0009	0.09%	99.91%
8	0.0003	0.03%	99.97%
9	0.0001	0.01%	99.99%
10	0.000045	0.005%	99.995%



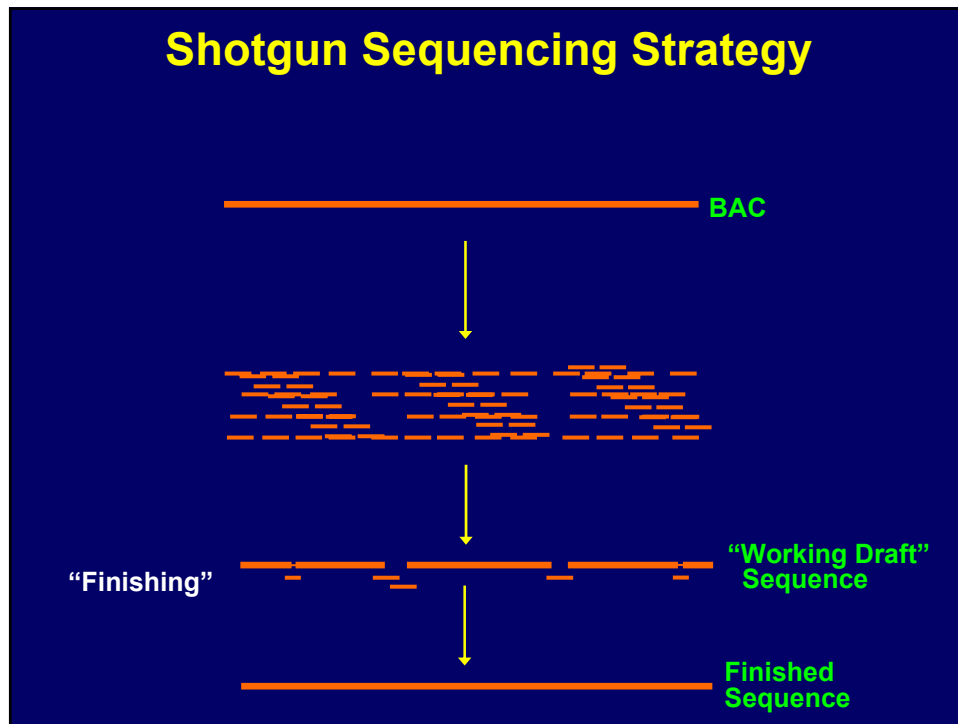
# Shotgun Sequence Assembly



“Consed” (Gordon et al., 1998)



## Shotgun Sequencing Strategy



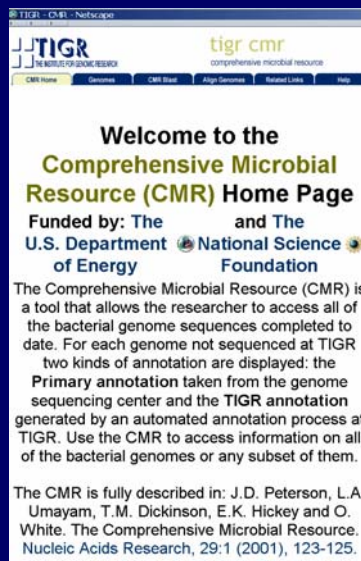
## Sequence Finishing: Resolving Ambiguities



\*\*\* Sequence Finishing: Remains Relatively Expensive \*\*\*

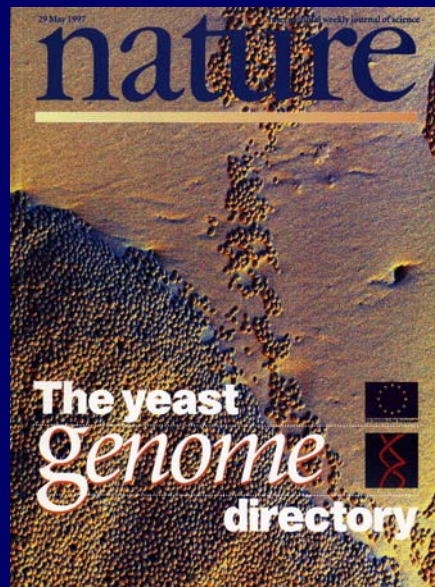
# Historically Significant Genome Sequencing Projects

## Bacterial Genome Sequences



[www.tigr.org](http://www.tigr.org)

## First Eukaryotic Genome Sequence



*Nature* 387:1-105, 1997

## First Animal Genome Sequence

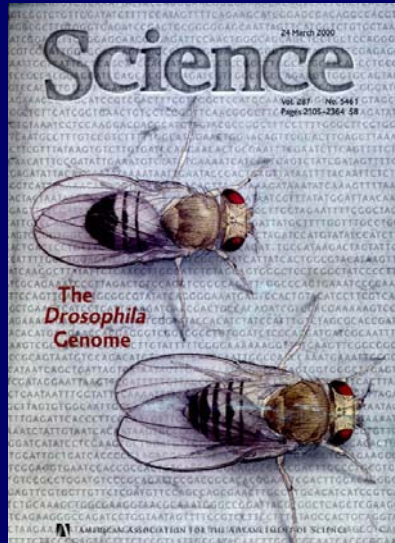


Genome Sequence of the Nematode *C. elegans*:  
A Platform for Investigating Biology

The *C. elegans* Sequencing Consortium\*

*Science* 282:1012-2018, 1998

## Second Animal Genome Sequence



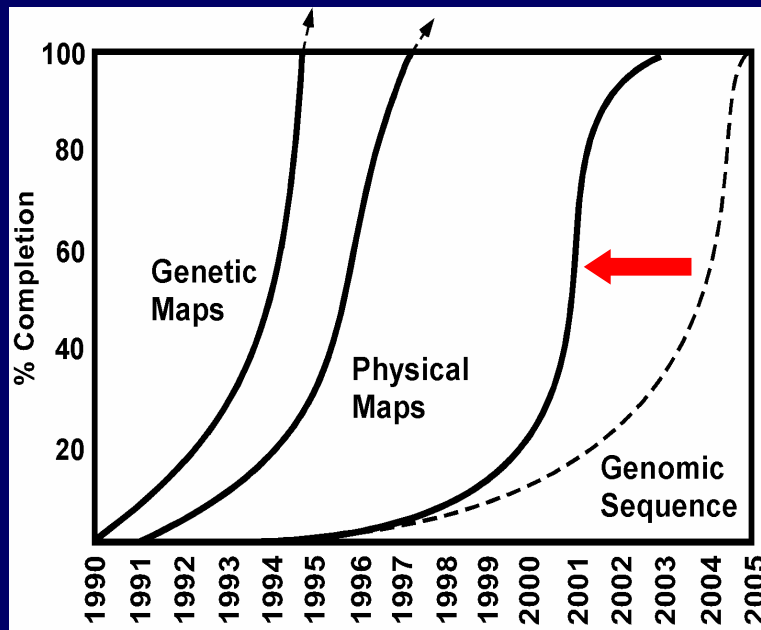
THE DROSOPHILA GENOME  
REVIEW

### The Genome Sequence of *Drosophila melanogaster*

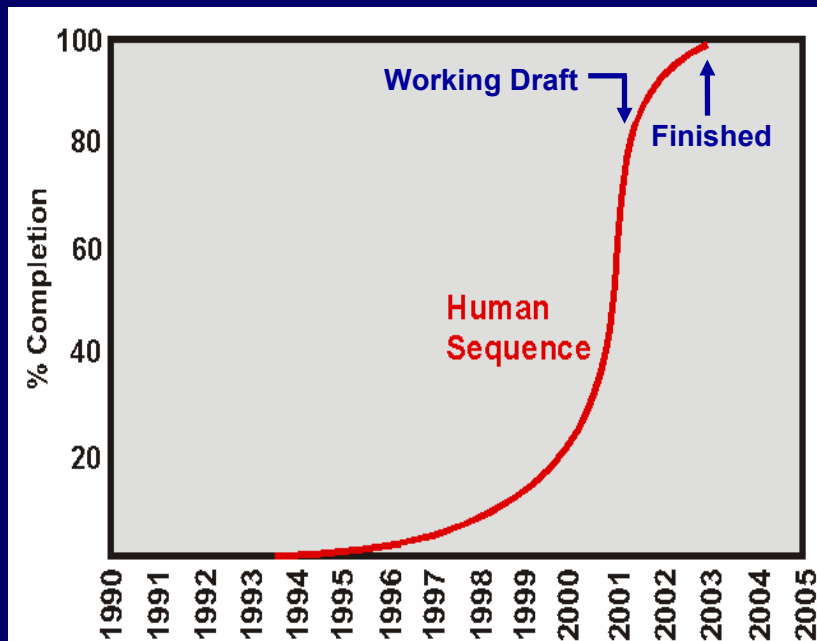
Mark D. Adams,<sup>1\*</sup> Susan E. Celniker,<sup>2</sup> Robert A. Holt,<sup>1</sup> Cheryl A. Evans,<sup>1</sup> Jeannine D. Gocayne,<sup>1</sup> Peter C. Amanatides,<sup>1</sup> Steven E. Scherer,<sup>1</sup> Peter W. Li,<sup>1</sup> Roger A. Hoskins,<sup>1</sup> Richard F. Gallie,<sup>1</sup> Read A. George,<sup>1</sup> Suzanna E. Lewis,<sup>1</sup> Stephen Richards,<sup>1</sup> Michael Ashburner,<sup>1</sup> Scott N. Henderson,<sup>1</sup> Granger G. Sutton,<sup>1</sup> Jennifer R. Wortman,<sup>1</sup> Mark D. Yandell,<sup>1</sup> Qing Zhang,<sup>1</sup> Lin X. Chen,<sup>1</sup> Rhonda C. Brandon,<sup>1</sup> Yu-Hui C. Rogers,<sup>1</sup> Robert G. Blazyn,<sup>1</sup> Mark Champe,<sup>1</sup> Barret D. Pfeiffer,<sup>1</sup> Kenneth H. Wan,<sup>1</sup> Clara Doyle,<sup>1</sup> Evan G. Baxter,<sup>1</sup> Gregg Helt,<sup>1</sup> Catherine R. Nelson,<sup>1</sup> George L. Gabor Miklos,<sup>1</sup> Joseph F. Abril,<sup>1</sup> Anna Aghayani,<sup>1</sup> Hai Jin An,<sup>1</sup> Cynthia Andrews-Pfannkoch,<sup>1</sup> Danita Baldwin,<sup>1</sup> Richard M. Ballew,<sup>1</sup> Anand Basu,<sup>1</sup> James Baxterdale,<sup>1</sup> Leyla Bayraktaroglu,<sup>1</sup> Eilan M. Beasley,<sup>1</sup> Karen Y. Besson,<sup>1</sup> P. V. Benos,<sup>1</sup> Benjamin P. Berman,<sup>1</sup> Deepali Bhandari,<sup>1</sup> Slave Bolshakov,<sup>11</sup> Dana Bokorova,<sup>12</sup> Michael R. Botchan,<sup>13</sup> John Book,<sup>14</sup> Peter Brokstein,<sup>15</sup> Philippe Brotille,<sup>14</sup> Kenneth C. Burtis,<sup>15</sup> Dana A. Besan,<sup>16</sup> Heather Butler,<sup>16</sup> Edouard Cadieu,<sup>17</sup> Angela Center,<sup>18</sup> Ishwar Chandra,<sup>19</sup> J. Michael Cherry,<sup>18</sup> Simon Cawley,<sup>18</sup> Carl Dahlke,<sup>18</sup> Lionel B. Davenport,<sup>19</sup> Peter Davies,<sup>19</sup> Beatriz de Pablos,<sup>20</sup> Arthur Delcher,<sup>21</sup> Zuoming Deng,<sup>21</sup> Anne Deslattes Hays,<sup>1</sup> Ian Dew,<sup>22</sup> Suzanne H. Dietz,<sup>23</sup> Kristina Dodson,<sup>24</sup> Lisa E. Doup,<sup>1</sup> Michael Dvornick,<sup>21</sup> Shannon Dugan-Rocha,<sup>1</sup> Boris C. Dunkov,<sup>25</sup> Patrick Dum,<sup>1</sup> Kenneth J. Durbin,<sup>2</sup> Carlos C. Evangelista,<sup>1</sup> Concepcion Ferraz,<sup>23</sup> Steven Ferreira,<sup>1</sup> Wolfgang Fleischmann,<sup>1</sup> Carl Foster,<sup>1</sup> Andrei E. Gabrielian,<sup>1</sup> Neha S. Garg,<sup>1</sup> William M. Galbert,<sup>1</sup> Ken Glasser,<sup>1</sup> Anna Glöckl,<sup>1</sup> Fangcheng Gong,<sup>1</sup> J. Harley Gorrell,<sup>26</sup> Zhiping Gu,<sup>1</sup> Ping Guan,<sup>1</sup> Michael Harris,<sup>1</sup> Naomi L. Harris,<sup>1</sup> Damon Harvey,<sup>1</sup> Thomas J. Heiman,<sup>1</sup> Judith K. Hernandez,<sup>1</sup> Jarrett Houck,<sup>1</sup> Damon Hostin,<sup>1</sup> Kathryn A. Houston,<sup>1</sup> Timothy J. Howland,<sup>1</sup> Ming-Hid Wei,<sup>1</sup> Chinyere Ibegwam,<sup>1</sup> Mena Jalali,<sup>1</sup> Francis Kalush,<sup>1</sup> Gary H. Karpen,<sup>27</sup> Zhaodl Ke,<sup>1</sup> James A. Kennison,<sup>28</sup> Karen A. Ketchum,<sup>1</sup> Bruce E. Kimmel,<sup>1</sup> Chinnappa D. Kodira,<sup>1</sup> Cheryl Kraft,<sup>1</sup> Soui Kravitz,<sup>1</sup> David Kulp,<sup>1</sup> Zhongqi Lu,<sup>1</sup> Paul Lasko,<sup>29</sup> Yiding Lei,<sup>1</sup> Alexander A. Levitsky,<sup>1</sup> Jianlin Li,<sup>1</sup> Zhanyu Li,<sup>1</sup> Yong Liang,<sup>1</sup> Xiaoying Lin,<sup>28</sup> Xiangjun Liu,<sup>1</sup> Bettina Mattel,<sup>1</sup> Tina C. McIntosh,<sup>1</sup> Michael P. McLeod,<sup>1</sup> Duncan McPherson,<sup>1</sup> Gennady Merkulov,<sup>1</sup> Natalia V. Milhina,<sup>1</sup> Clark Moberly,<sup>1</sup> Joe Morris,<sup>1</sup> All Hoshreli,<sup>1</sup> Stephen H. Mount,<sup>27</sup> Mei Moy,<sup>1</sup> Brian Murphy,<sup>1</sup> Lee Murphy,<sup>28</sup> Donna M. Murray,<sup>1</sup> David L. Nelson,<sup>29</sup> David R. Nelson,<sup>29</sup> Keith A. Nelson,<sup>1</sup> Katherine Nilson,<sup>1</sup> Deborah R. Nusser,<sup>1</sup> Joanne M. Paclab,<sup>1</sup> Michael Palazzolo,<sup>2</sup> Gjang S. Pittman,<sup>1</sup> Sue Pan,<sup>1</sup> John Pollard,<sup>1</sup> Vinita Puri,<sup>1</sup> Martin G. Reese,<sup>1</sup> Knut Reinert,<sup>1</sup> Karin Remington,<sup>1</sup> Robert D. C. Saunders,<sup>28</sup> Frederick Scheeler,<sup>1</sup> Hua Shen,<sup>1</sup> Biliang Christopher Shou,<sup>1</sup> Inga Sidén-Kiamos,<sup>1</sup> Michael Simpson,<sup>1</sup> Marian P. Skupski,<sup>1</sup> Tom Smith,<sup>1</sup> Eugene Spler,<sup>1</sup> Allan C. Spreading,<sup>1</sup> Mark Stapleton,<sup>2</sup> Renee Strong,<sup>1</sup> Eric Sun,<sup>1</sup> Robert Svirskii,<sup>22</sup> Cyndee Tector,<sup>1</sup> Russell Turner,<sup>1</sup> Eli Venter,<sup>1</sup> Alhui H. Wang,<sup>1</sup> Xin Wang,<sup>1</sup> Zhen-Yuan Wang,<sup>1</sup> David A. Wassarman,<sup>23</sup> George H. Weinstock,<sup>1</sup> Jean Weissenbach,<sup>1</sup> Sherita M. Williams,<sup>1</sup> Trevor Woodruff,<sup>1</sup> Kim C. Worley,<sup>1</sup> David Wu,<sup>1</sup> Song Yang,<sup>1</sup> Q. Allison Yao,<sup>1</sup> Jian Ye,<sup>1</sup> Ru-Fang Yeh,<sup>1</sup> Jaystree S. Zaveri,<sup>1</sup> Ming Zhang,<sup>1</sup> Guangren Zhang,<sup>1</sup> Qi Zhao,<sup>1</sup> Lianheng Zheng,<sup>1</sup> Xiangqun H. Zheng,<sup>1</sup> Fei N. Zhong,<sup>1</sup> Wenyan Zhong,<sup>1</sup> Xiaojun Zhou,<sup>1</sup> Shaoping Zhu,<sup>1</sup> Xiaohong Zhu,<sup>1</sup> Hamilton O. Smith,<sup>1</sup> Richard A. Gibbs,<sup>1</sup> Eugene W. Myers,<sup>1</sup> Gerald M. Rubin,<sup>24</sup> J. Craig Venter<sup>1</sup>

Science 287:2185-2195, 2000

## Revised Timetable for Human Genome Sequencing



### Timetable for Human Genome Sequencing



### Human Genome Sequencing Centers



Whitehead Institute/MIT  
Genome Sequencing Center



JGI  
JOINT GENOME INSTITUTE



## Human Genome Sequencing Centers



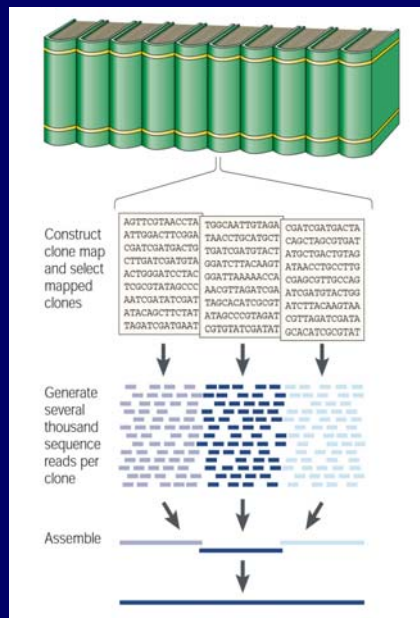
## June, 2000 Announcement



## February, 2001 Publications



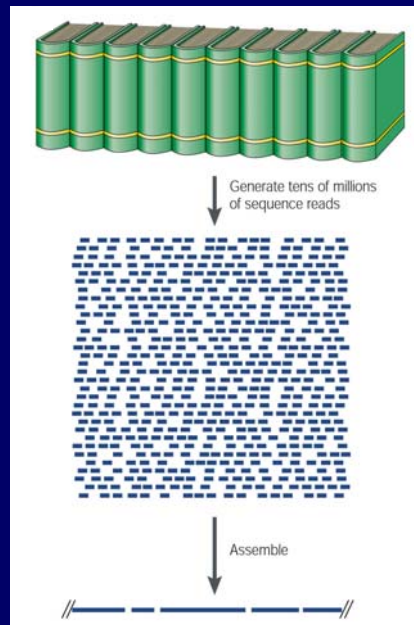
## BAC-by-BAC Shotgun Sequencing



Green (2001)

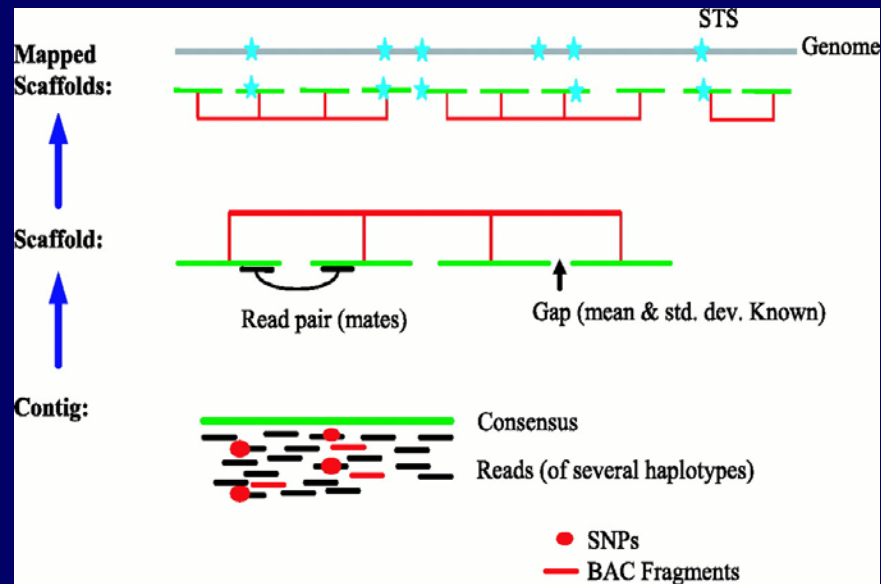


## Whole-Genome Shotgun Sequencing



Green (2001)

## Whole-Genome Shotgun Sequence Assembly



Venter et al., 2001

## April, 2003 Completion



## International Human Genome Sequencing Consortium



- 6 Countries
- 20 Sequencing Centers
- 1000's of Individuals
- ~1,000 bases per second, 24 hours per day, 7 days per week



108TH CONGRESS  
1ST SESSION

## S. CON. RES. 10

Designating April 2003 as "Human Genome Month" and April 25 as "DNA Day".

---

IN THE SENATE OF THE UNITED STATES

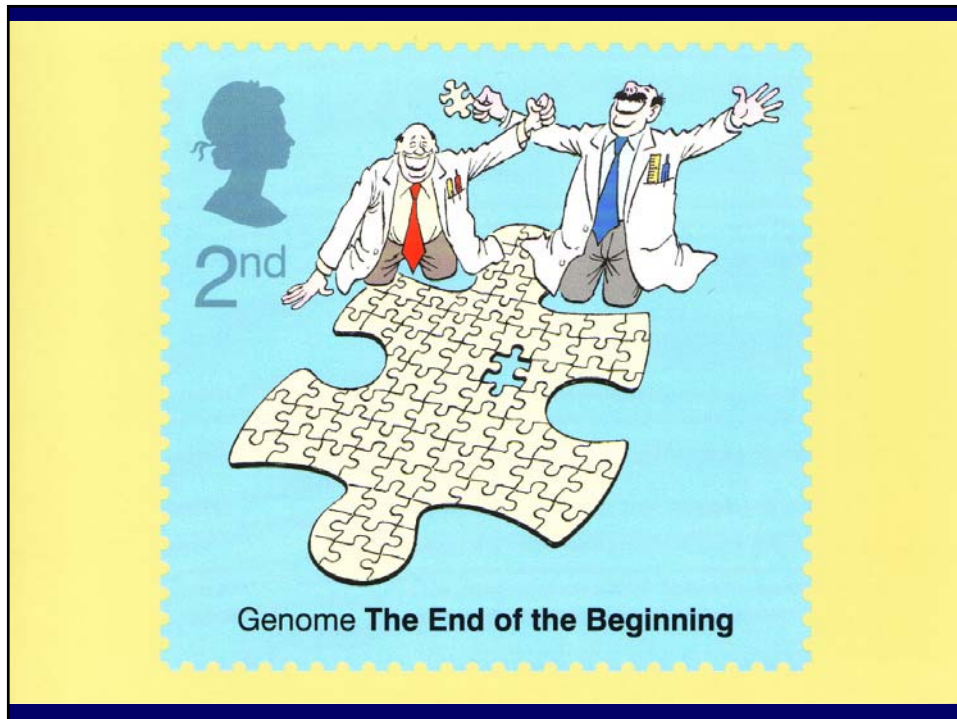
FEBRUARY 27, 2003

Mr. GREGG (for himself, Mr. KENNEDY, Ms. SNOWE, and Mr. DASCHLE) submitted the following concurrent resolution; which was considered and agreed to

---

### CONCURRENT RESOLUTION

Designating April 2003 as "Human Genome Month" and April 25 as "DNA Day".



**October, 2004 Publication**

**Tetraodon to human**  
Evolutionary history in genome sequences

**General relativity**  
Did the orbit move for you?

**The human genome**  
Going the last mile

**Antibiotics crisis**  
Market forces fail to deliver

**Medical ethics**  
Choosing deafness

**articles**

**Finishing the euchromatic sequence of the human genome**

*International Human Genome Sequencing Consortium\**

\* A list of authors and their affiliations appears in the Supplementary Information.

The sequence of the human genome encodes the genetic instructions for human physiology, as well as rich information about human evolution. In 2001, the International Human Genome Sequencing Consortium reported a draft sequence of the euchromatic portion of the human genome. Since then, the international collaboration has worked to convert this draft into a genome sequence with high accuracy and nearly complete coverage. Here, we report the result of this finishing process. The current genome sequence (Build 25) contains 2.9 billion nucleotides interrupted by only 241 gaps. Coverage – 99% of the euchromatic genome and is accurate to an error rate of ~1 error per 100,000 bases. Many of the remaining euchromatic gaps are associated with segmental duplications and will require focused work with new methods. The non-coding sequence, the first for a vertebrate, greatly improves the precision of biological analyses of the human genome including studies of gene number, birth and death. Notably, the human genome encodes only 28,000–31,000 protein-coding genes. The genome sequence reported here should serve as a firm foundation for biomedical research in the decades ahead.

The Human Genome Project (HGP) was launched in 1990 with the goal of obtaining a high-quality sequence of the non-repetitive portion of the human genome. The initial work followed a two-pronged approach: (1) the mapping of the human and mouse genomes<sup>1,2</sup> to allow the study of inherited disease and provide a scaffold for genome assembly; and (2) the sequencing of organisms with smaller, single genomes<sup>3,4</sup> to serve as model for method development and practice in sequencing the human genome. With success along both paths, the sequencing of the human genome itself eventually became feasible. The International Human Genome Sequencing Consortium (IHGSC), an open collaboration involving twenty centres in six countries, was formed to carry out this component of the HGP.

In February 2001, the IHGSC<sup>5</sup> and Celera Genomics<sup>6</sup> each reported draft sequence genomes providing a first overall view of the human genome. These sequences allowed systematic study of the human genome itself, including identification of genes, functional annotation of proteins, regional differences in genome composition, distribution and history of transposable elements, distribution of polymorphisms, and relationships between genetic susceptibility and physical disease. Moreover, systematic knowledge of the human genome has enabled new methods of approach that have markedly accelerated subsequent developments. The IHGSC sequenced an average of ~10% of the euchromatic genome. It was interrupted by ~150,000 gaps and the order and orientation of many sequences could not be assigned but were not established. The IHGSC then turned to the challenge of completing the euchromatic genome. Operationally, a finished sequence was defined as having an error rate of, at most, one error per 10<sup>5</sup> bases, and the final, complete coverage of a finished sequence of at least 95% of the euchromatic genome, with the only gaps being those associated with segmental duplications<sup>7</sup> (see <http://www.genome.gov/109992>). The goal was challenging because the human genome is complex with much repetitive content and large segmental duplications, which greatly complicate the determination of precise sequence and repeat length. The most complete sequences have been obtained so far with the three nucleotide sequencing technologies<sup>8,9</sup>, repeated work<sup>10</sup> and the 'rafting'<sup>11</sup>. These processes are all roughly 10-fold slower than the human genome and have much smaller coverage.

We describe here the results of a major effort by the IHGSC towards the goal of a complete human sequence. The number of gaps that have reduced 80-fold to only 161, most of which are associated with segmental duplications and will require new methods for resolution. The assembled non-coding genome sequence has an error rate of only ~1 error per 100,000 bases. It contains 241 gaps, nucleotide and gaps – 99% of the euchromatic genome. The paper describes the current genome sequence and the process used to produce it, examines the accuracy and completeness of the sequence, and illustrates biological features made possible by the sequence. We also attempt here a comparative analysis of the contents of the human genome. An initial analysis was previously reported<sup>12</sup> and a series of papers is being written describing the individual chromosomes<sup>13–15</sup>, including association of genes and other features.

**Current genome sequence**

**Finishing process**

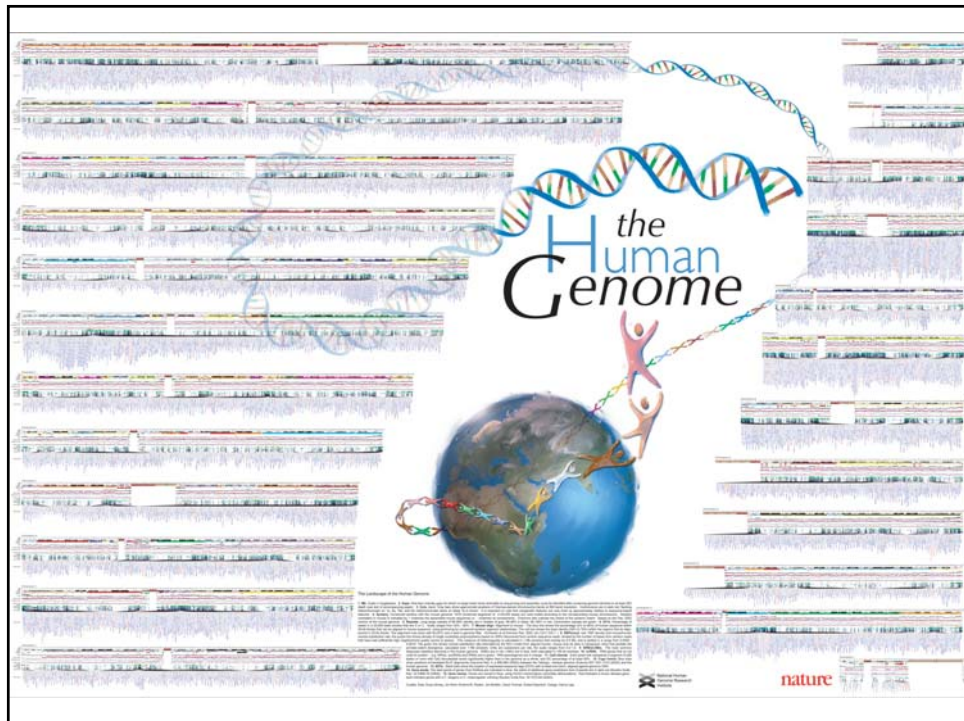
The process of converting the initial draft sequence into a more complete sequence is referred to as 'finishing'. It is a complex iterative process that proceeds simultaneously at multiple scales, ranging from single nucleotides to the lengths of whole chromosomes. The fundamental challenge that this process requires that are not well represented in draft is that through random changes sequencing used to be highly enriched for problematic sequences. Finishing such regions required the development of special approaches, which varied substantially over time and varied among centres.

Initially, the finishing process involved two distinct components: (1) producing finished maps, consisting of continuous and clone maps of overlapping large insert clones spanning the euchromatic region of chromosomes; and (2) producing finished clones, consisting of continuous and accurate nucleotide sequence across each large insert clone. In practice, these two components were tightly interrelated in that progress in each often depended on results from the other. The sequence information that was most highly interested in that progress in each often depended on results from the other. The sequence information that was most highly interested in that progress in each often depended on results from the other. The sequence information that was most highly interested in that progress in each often depended on results from the other.

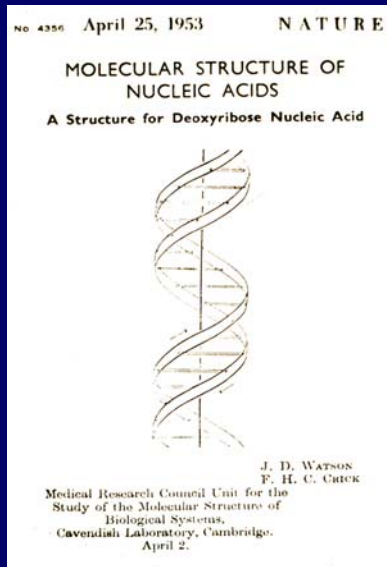
The final, non-generated sequence comprises 2,911,319,271 base-pair human DNA (total length ~3.14 gigabases (Gb)) and finished the sequence from 45.7% of their gene total length, 1,672.0 Gb. The clones contained primarily of bacterial artificial chromosomes (BACs).

© 2004 Nature Publishing Group

Nature 431:931-945, 2004

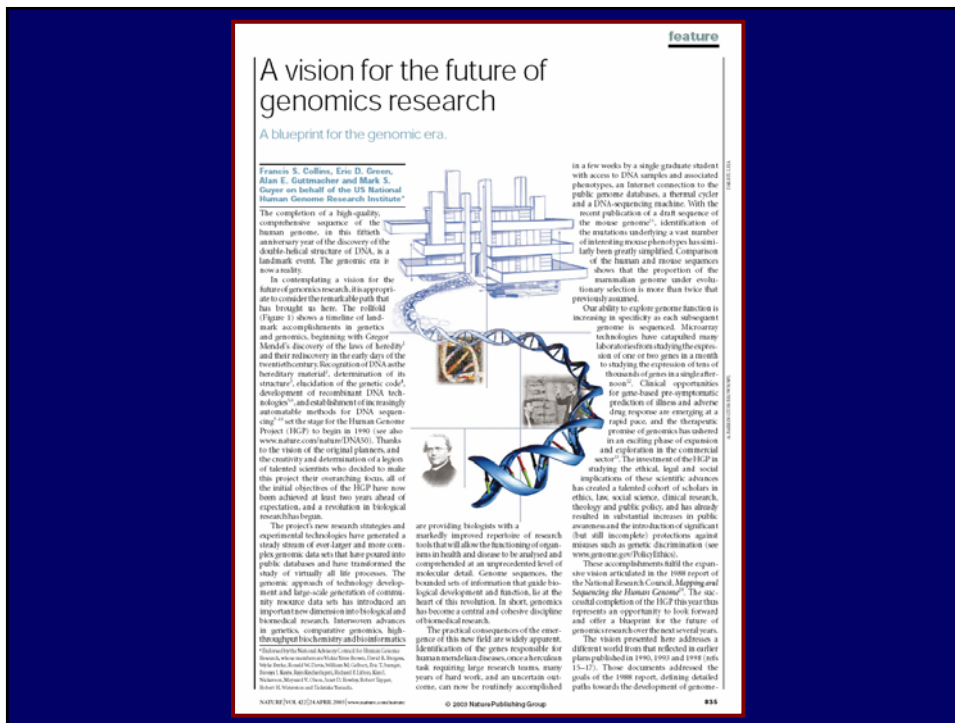


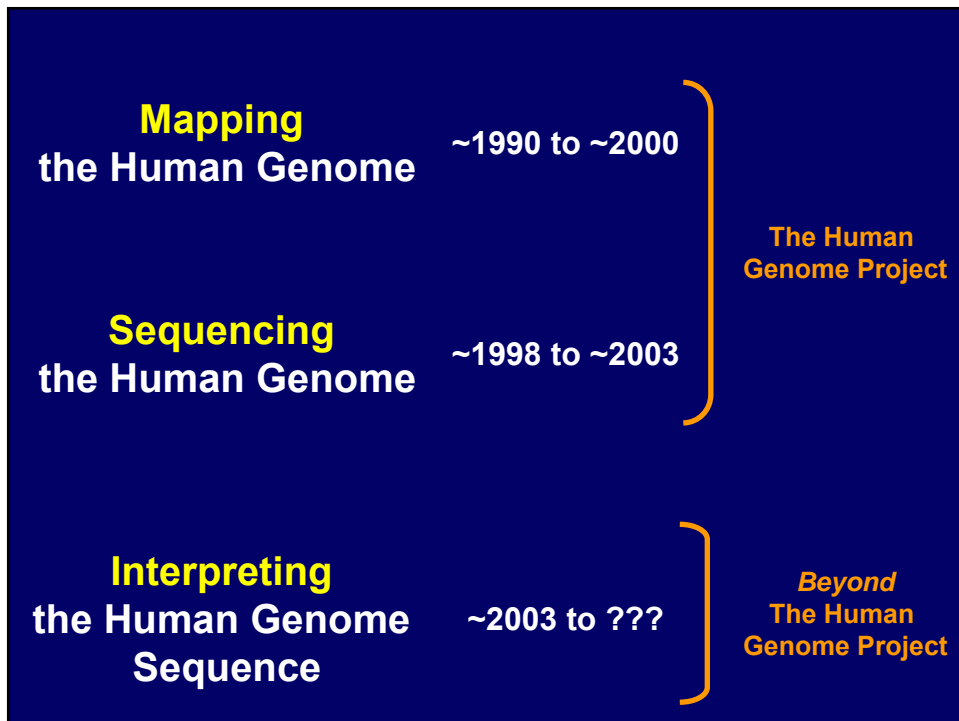
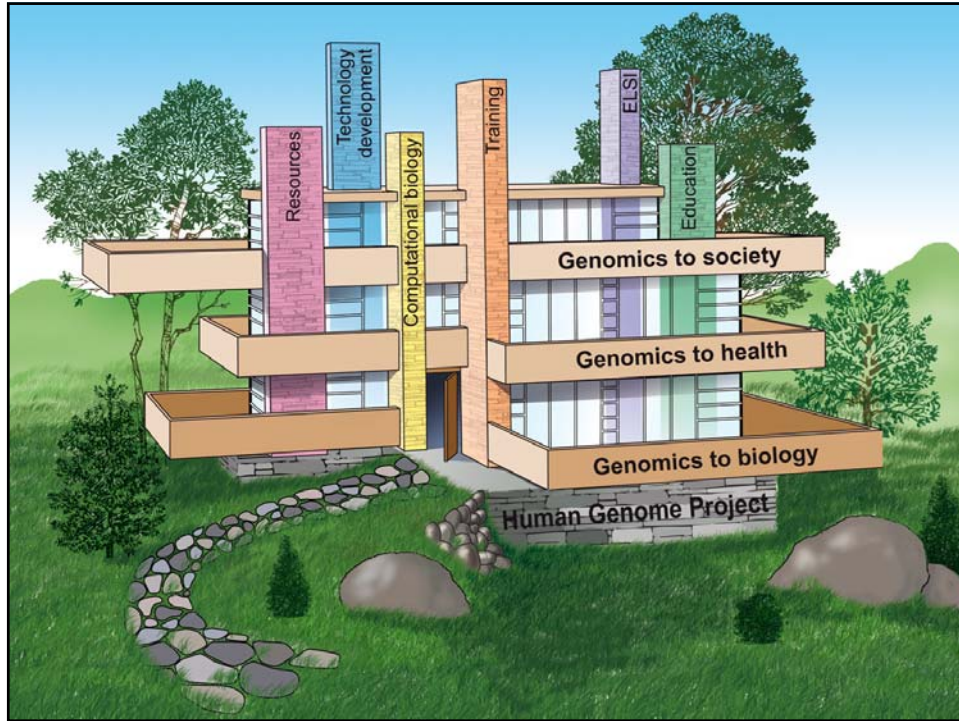
**April, 1953** → **April, 2003**



**All of the original goals of the  
Human Genome Project have  
been accomplished!**

**What's Next?**









## Comparative Sequence Analysis

*Using the Experiments of Evolution to Decode the Human Genome*

**Species A**

```

GATCGTCTAGAACTCGAGATC
TCTGAGATCTCTGGAAACTGT
GTGATGTGACTAGCCACAGTTA
GTTACGTGTGAGAGATGTATGA
TGCACCTGACCCGGGTTCACT
CTCAAGACTCACTCCACTCA
GAGGCCACCCGCTGTGCAC
TACCGATATACGATACCTAC
ACAGGTGTGACACACCCTACC
CSTCCACACAGACTCACTCC
ACCCTCAGAGGCCACCCGCGCT
GTGCACTACCGATACAGCAT
ACCACACAGGTGTGACACAGT
ATCCTTACACACTTACACATT
ACCATATATCCACTACACAC
ATACTTACCCCTTGCACACT
ATTATTATACCGGACCGGAG
                    
```

**Compare**

**Species B**

```

TATCGGCTAGAACTCGAGATC
TCTGAGATCTCTGGAAACTGT
GTGATGTGACTAGCCACAGTTA
GTTGAGAGATGTATGATGCA
CTTGACCCGGGTTCACTCTCA
ACGACTCACTCCACTCAGAGG
CCACCCGCTGTGCACAGTCC
ACCATGATCTTACACACTTA
CACATCACTCTCAAGACTCAC
TCCACTCAGAGGCCACCGCC
GCTGTGACGTCCACACAGATC
CTTACACACTTACACATTACC
ATATATCCACTACACACATA
CCTTACCATATATCCACTACC
ACCATACTACCCCTTGCAC
ACCTATTATTATACCGGGA
GAGGGTGCACACTGTGACA
                    
```

```

GATCGTCTAGAACTCGAGATC
TCTGAGATCTCTGGAAACTGT
GTGATGTGACTAGCCACAGTTA
CCTGACCCGGGTTCACTCTCA
ACGACTCACTCCACTCAGAGG
CCACCCGCTGTGCACAGTCC
ACCATGATCTTACACACTTA
CACATCACTCTCAAGACTCAC
TCCACTCAGAGGCCACCGCC
GCTGTGACGTCCACACAGATC
CTTACACACTTACACATTACC
ATATATCCACTACACACATA
CCTTACCATATATCCACTACC
ACCATACTACCCCTTGCAC
ACCTATTATTATACCGGGA
GAGGGTGCACACTGTGACA
                    
```

Sequences in Common (i.e., 'Conserved')

## Comparing Genomes is Like Cryptography

```

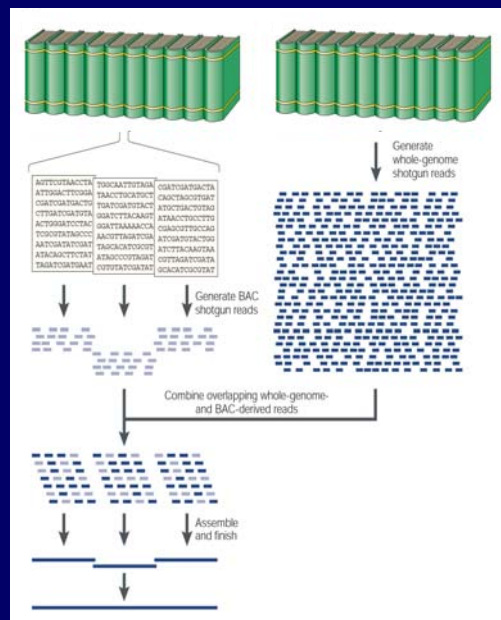
CKQEBHEREYTWASULSCZMEISDFOGETHEBLPBGODFQSTLKSUFFRTAC
DLUCEHEREZBRTTOISAWNDCDARJJPThERROFGODERGHCLSUFFBRHA
                    
```

## Functional Elements: Coding vs. Non-Coding

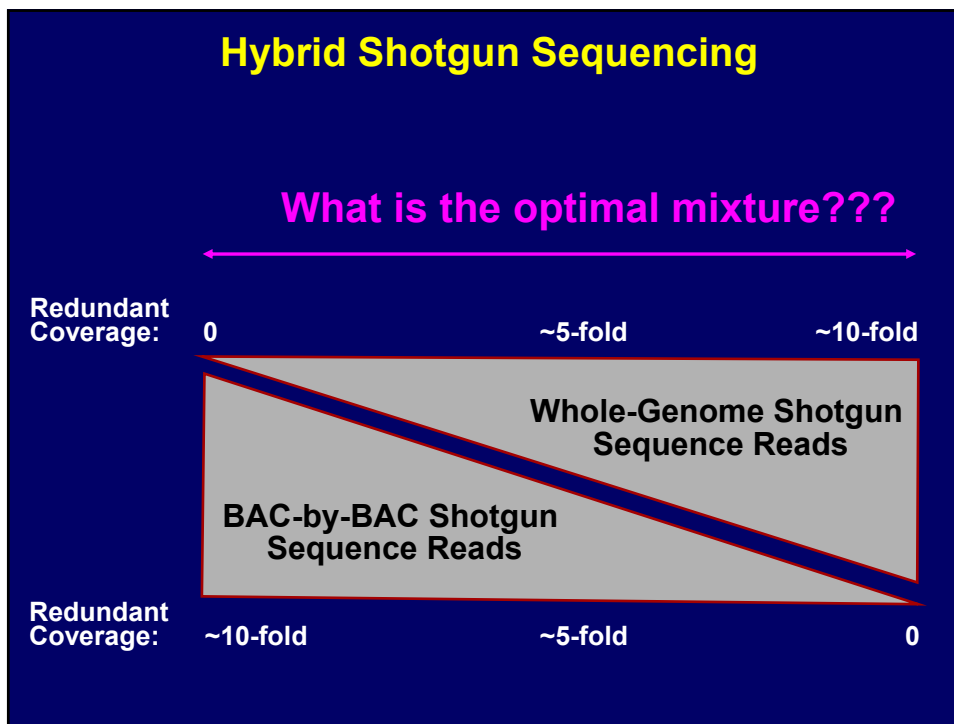
- **Coding Sequences (i.e., Genes)**
  - Relatively EASY to Identify
  - Mostly Know What to Look For
  - Complementary Data Sets Available (ESTs, cDNAs)
  - Ever-Improving Computational Gene Predictions
- **Non-Coding Functional Sequences**
  - HARD to Identify
  - Very Little Known About What to Look For
  - Virtually No Complementary Data Sets Available
  - Poor Computational Predictions

**Major role for comparative sequence analysis will be the identification of functionally important, non-coding sequences**

## Hybrid Shotgun Sequencing



Green (2001)



**nature**  
December 2002  
International weekly journal of science

**The mouse genome**  
Experimental model for human biology

Atmospheric CO<sub>2</sub>: A drop in the ocean  
Drug reduction: Flexibility in trials  
Regulatory T cells: Basis for persistent infection

**Nature 420:520-562, 2002**

**articles**  
**Initial sequencing and comparative analysis of the mouse genome**  
Mouse Genome Sequencing Consortium\*

**nature**  
1 April 2004  
International weekly journal of science

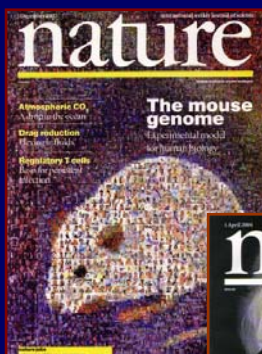
**The rat genome**  
Insights into mammalian evolution

Superconductivity: Diamond springs a surprise  
Office life: Make those e-mails count  
SARS vaccine: Immunity induced in mice

**Nature 428:493-521, 2004**

**articles**  
**Genome sequence of the Brown Norway rat yields insights into mammalian evolution**  
Rat Genome Sequencing Project Consortium\*

## Human-Rodent Sequence Comparisons



Nature  
420:520-562, 2002



Nature  
428:493-521, 2004

- ~40% in Alignments
- ~5% Under Selection
- ~1.5% Protein Coding
- ~3.5% Non-Coding

## Multi-Species Sequence Comparisons

GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG	GATCCTCTAGAACTCTCG AGATCTCTGAGAGTCCTG GGGAAACTCTGTGTGATGT GACTTCCACAGCTTACCG TGTGAGAGATGTATGAT GCACCTGACCCCGGTTT GACTCTCAACAGCTTACCG TCCACTCTGAGAGTCCCA CCGCGCTGTGACAGTCC CACCAGATCTCTTACCA CAGCTTCAATTAACAG ATACTCCACTACCCAGAC ATACTCCACTACCCAGAC CAGCTTATATATATACCG
--	--	--	--	--	--	--

HUMAN

## Multi-Species Comparative Sequence Analysis

### Comparative analyses of multi-species sequences from targeted genomic regions

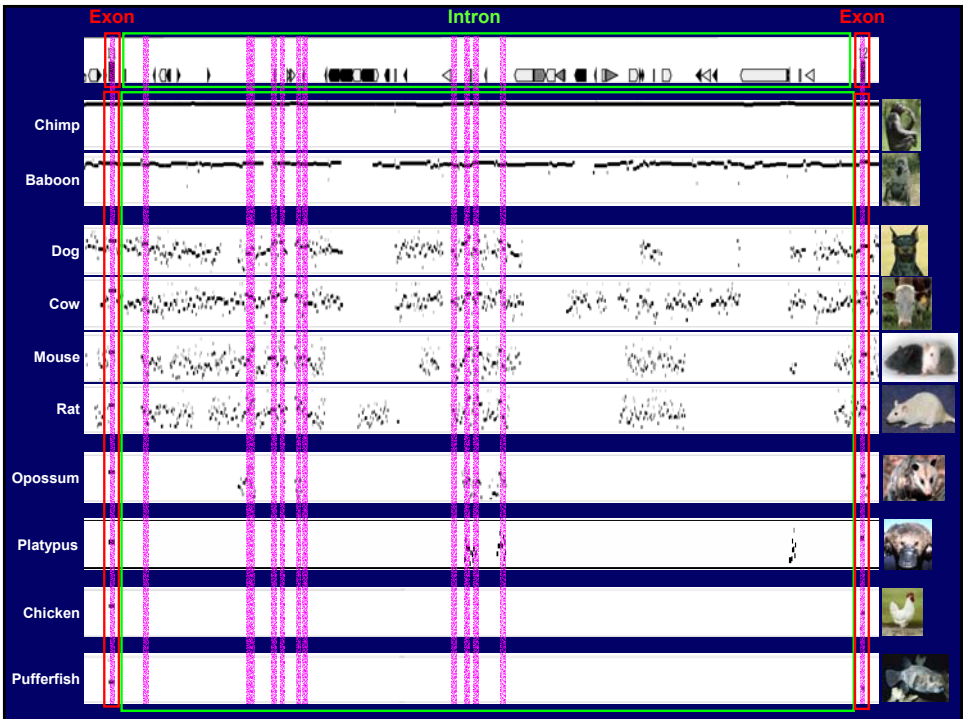
J. W. Thomas<sup>1</sup>, J. W. Tomblin<sup>2</sup>, A. W. Rabeck<sup>3</sup>, A. G. Buehler<sup>4</sup>, S. M. Baskin-Shribler<sup>5</sup>, E. S. Bergler<sup>6</sup>, M. Brinkley<sup>7</sup>, S. C. Cooper<sup>8</sup>, P. A. D'Amico<sup>9</sup>, G. Hochstadt<sup>10</sup>, R. Minkley<sup>11</sup>, W. G. Rouse<sup>12</sup>, M. S. Scheraga<sup>13</sup>, S. J. Singer<sup>14</sup>, W. J. Wolf<sup>15</sup>, D. Zerkow<sup>16</sup>, T. G. Brown<sup>17</sup>, R. Dwyer<sup>18</sup>, B. E. Galley<sup>19</sup>, S. Hahn<sup>20</sup>, L. Hahn<sup>21</sup>, J. S. Hall<sup>22</sup>, A. H. Patel<sup>23</sup>, S.-H. Lee<sup>24</sup>, V. S. R. Muth<sup>25</sup>, S. J. Somers<sup>26</sup>, R. S. Paranjape<sup>27</sup>, S. L. Scherer<sup>28</sup>, M. Miller<sup>29</sup>, A. Agre<sup>30</sup>, S. Stroganov<sup>31</sup>, R. Corbett<sup>32</sup>, C. P. Brinkley<sup>33</sup>, S. T. Shultz<sup>34</sup>, S. Osofsky<sup>35</sup>, L. Olson<sup>36</sup>, J. Gupta<sup>37</sup>, P. Hahnle<sup>38</sup>, J. C. Lee<sup>39</sup>, M. C. Hwang<sup>40</sup>, J. Kasper<sup>41</sup>, J. L. Grier<sup>42</sup>, R. Leger<sup>43</sup>, M. J. Lee<sup>44</sup>, G. L. Muth<sup>45</sup>, C. A. Maudsley<sup>46</sup>, S. S. Maudsley<sup>47</sup>, J. C. McInerney<sup>48</sup>, R. Pearce<sup>49</sup>, J. Shultz<sup>50</sup>, E. S. Thompson<sup>51</sup>, J. T. Drenth<sup>52</sup>, C. Thompson<sup>53</sup>, J. L. Vogel<sup>54</sup>, M. A. Webster<sup>55</sup>, A. S. Webster<sup>56</sup>, J. S. Wagner<sup>57</sup>, A. C. Wang<sup>58</sup>, L.-M. Zhang<sup>59</sup>, K. Ouyang<sup>60</sup>, R. Zhu<sup>61</sup>, R. Zhu<sup>62</sup>, C. L. Shih<sup>63</sup>, P. J. de Jong<sup>64</sup>, G. S. Lawrence<sup>65</sup>, A. P. Tang<sup>66</sup>, A. Chakravarti<sup>67</sup>, D. Haussler<sup>68</sup>, P. Green<sup>69</sup>, M. Miller<sup>70</sup> & E. S. Green<sup>71</sup>

<sup>1</sup>Genome Technology Branch, National Human Genome Research Institute, and NIH Intron Sequencing Center, National Institutes of Health, Bethesda, Maryland 20892, USA  
<sup>2</sup>Center for Biomolecular Science and Engineering, University of California, Santa Cruz, California 95064, USA  
<sup>3</sup>Department of Genetics, Johns Hopkins University School of Medicine, Baltimore, Maryland 21207, USA  
<sup>4</sup>Department of Computer Science and Engineering, The Pennsylvania State University, University Park, Pennsylvania 16802, USA  
<sup>5</sup>Children's Hospital Oakland Research Institute, Oakland, California 94612, USA  
<sup>6</sup>The Weizmann Center for Laboratory and Research, New York State Department of Health, Albany, New York 12242, USA  
<sup>7</sup>The Institute for System Biology, Seattle, Washington 98195, USA  
<sup>8</sup>Medical Physics Medical Institute, University of California, Santa Cruz, California 95064, USA  
<sup>9</sup>Howard Hughes Medical Institute and Department of Genome Sciences, University of Washington, Seattle, Washington 98195, USA

<sup>10</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>11</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>12</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>13</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>14</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>15</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>16</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>17</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>18</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>19</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>20</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>21</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>22</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>23</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>24</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>25</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>26</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>27</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>28</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>29</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>30</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>31</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>32</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>33</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>34</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>35</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>36</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>37</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>38</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>39</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>40</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>41</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>42</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>43</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>44</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>45</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>46</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>47</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>48</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>49</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>50</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>51</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>52</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>53</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>54</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>55</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>56</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>57</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>58</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>59</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>60</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>61</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>62</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>63</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>64</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>65</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>66</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>67</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>68</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>69</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>70</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA  
<sup>71</sup>Present address: Department of Human Genetics, Harvard Medical School, Boston, MA 02115, USA

- Targeted Genomic Regions
- BAC-Based Sequencing in Multiple Vertebrates
- Identify Highly Conserved Non-Coding Sequences
- Conserved Sequences Correlate with Functional Elements

Thomas et al. (2003)



## Additional Vertebrate Genome Sequencing Efforts



Chimpanzee



Macaque



Dog



Cow



Monodelphis



Chicken



Xenopus

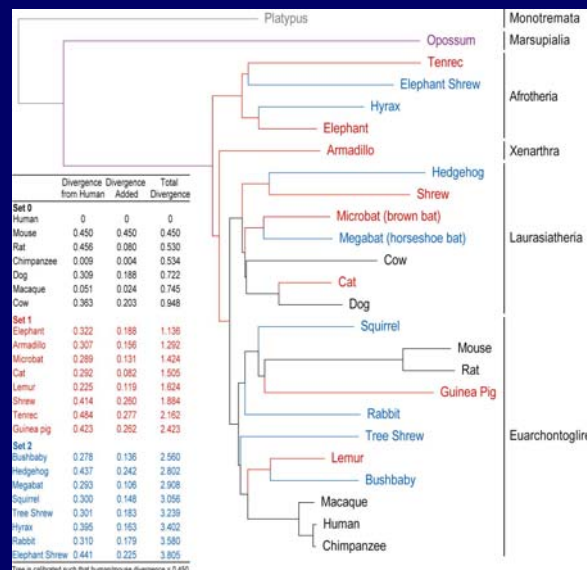


Zebrafish

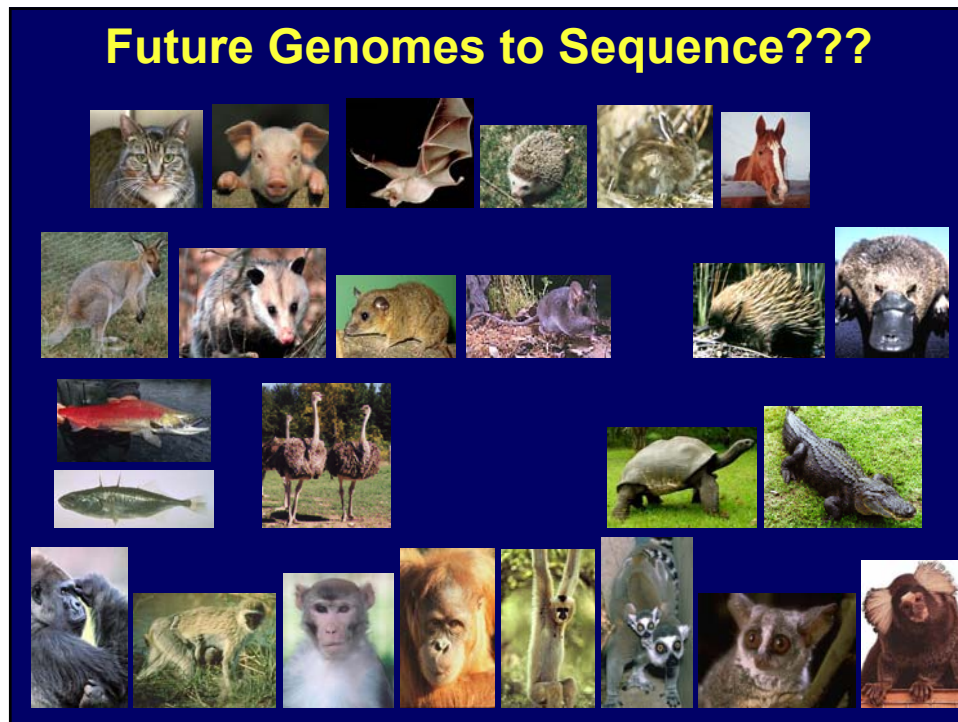


Pufferfish

## Low-Redundancy Sequencing of Multiple Vertebrate Genomes



Margulies et al., *PNAS*, in press, 2005



## ENCODE Project

- **ENCODE: ENCyclopedia Of DNA Elements**
- **Goal: Compile a *comprehensive encyclopedia* of all functional elements in the human genome**
- **Initial pilot project: 1% of human genome**
- **Apply multiple approaches to study and analyze that 1% in a consortium fashion**



## ENCODE Project: Web Sites

The left screenshot displays the ENCODE Project website on genome.gov. It includes a navigation menu and a main section titled "The ENCODE Project: Encyclopedia of DNA Elements". Below this, there is a section for "ENCODE Target Regions (January 2004)" which contains a table of genomic regions. The table has columns for "Region", "Description", and "Chr". The right screenshot shows the UCSC Genome Browser interface on genome.ucsc.edu, displaying a genomic track for chromosome 11. The track shows various genomic features and tracks, including a track for "ENCODE Target Regions".

[genome.gov/ENCODE](http://genome.gov/ENCODE)

[genome.ucsc.edu/ENCODE](http://genome.ucsc.edu/ENCODE)

## Current Big Challenges...

- Defining “Saturation Points” in Terms of Information Gained by Comparative Sequence Analyses
- The “\$1000 Genome”
- Medical Sequencing (aka, Human Re-Sequencing)

## Bibliography

Adams, M.D. et al. (2000). The genome sequence of *Drosophila melanogaster*. *Science* 287, 2185-2195.

Aparicio, S. et al. (2002). Whole-genome shotgun assembly and analysis of the genome of *Fugu rubripes*. *Science* 297, 1301-1310.

Birren, B. et al. (1998). Bacterial artificial chromosomes. In *Genome Analysis: A Laboratory Manual, Vol. 3 Cloning systems* (B. Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 241-295.

*C. elegans* Sequencing Consortium (1998). Genome sequence of the nematode *C. elegans*: a platform for investigating biology. *Science* 282, 2012-2018.

Collins, F.S. et al. (2003). A vision for the future of genomics research: a blueprint for the genomic era. *Nature* 422, 835-847.

Rat Genome Sequencing Project Consortium (2004). Genome sequence of the Brown Norway rat yields insights into mammalian evolution. *Nature* 428, 493-521.

Goffeau, A. et al. (1997). The Yeast Genome Directory. *Nature* 387S, 1-105.

Gordon, D. et al. (1998). Consed: a graphical tool for sequence finishing. *Genome Research* 8, 195-202.

Green, E.D. (2001). Strategies for the systematic sequencing of complex genomes. *Nature Rev Genet* 2, 573-583.

Green, E.D. (2001). The Human Genome Project and its impact on the study of human disease. In *The Metabolic and Molecular Bases of Inherited Disease, 8th Edition* (C.R. Scriver, A.L. Beaudet, W.S. Sly, D. Valle, B. Childs, K.W. Kinzler, and B. Vogelstein, eds.; McGraw-Hill, Inc.), pp. 259-298.

Green, E.D. et al. (1998). Yeast artificial chromosomes. In *Genome Analysis: A Laboratory Manual, Vol. 3 Cloning systems* (B. Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 297-565.

Gregory, S.G. et al. (1997). Genome mapping by fluorescent fingerprinting. *Genome Research* 7, 1162-1168.

Hillier, L.W. et al. (2004). Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* 432, 695-716.

International Human Genome Sequencing Consortium (2001). Initial sequencing and analysis of the human genome. *Nature* 409, 860-921.

International Human Genome Sequencing Consortium (2004). Finishing the euchromatic sequence of the human genome. *Nature* 431, 931-945.

Jaillon, O. et al. (2004). Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. *Nature* 431, 946-957.

Margulies, E.M. et al. (2005). An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing. *Proc Natl Acad Sci*, in press.

Marra, M.A. et al. (1997). High throughput fingerprint analysis of large-insert clones. *Genome Research* 7, 1072-1084.

Messing, J. and Llaca, V. (1998). Importance of anchor genomes for any plant genome project. *Proc Natl Acad Sci* 95, 2017-2020.

Miller, W. et al. (2004). Comparative genomics. *Ann Rev Genomics Hum Genet* 5, 15-56.

Mouse Genome Sequencing Consortium (2002). Initial sequencing and comparative analysis of the mouse genome. *Nature* 420, 520-562.

Shizuya, H. et al. (1992). Cloning and stable maintenance of 300-kilobase-pair fragments of human DNA in *Escherichia coli* using an F-factor-based vector. *Proc Natl Acad Sci* 89, 8794-8797.

Thomas, J.W. et al. (2003). Comparative analyses of multi-species sequences from targeted genomic regions. *Nature* 424, 788-793.

Venter, J.C. et al. (2001). The sequence of the human genome. *Science* 291, 1304-1351.

Wilson, R.K. and Mardis, E.R. (1997). Fluorescence-based DNA sequencing. In *Genome Analysis: A Laboratory Manual, Vol. 1 Analyzing DNA* (B. Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 301-395.

Wilson, R.K. and Mardis, E.R. (1997). Shotgun sequencing. In *Genome Analysis: A Laboratory Manual, Vol. 1 Analyzing DNA* (B. Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 397-454.